

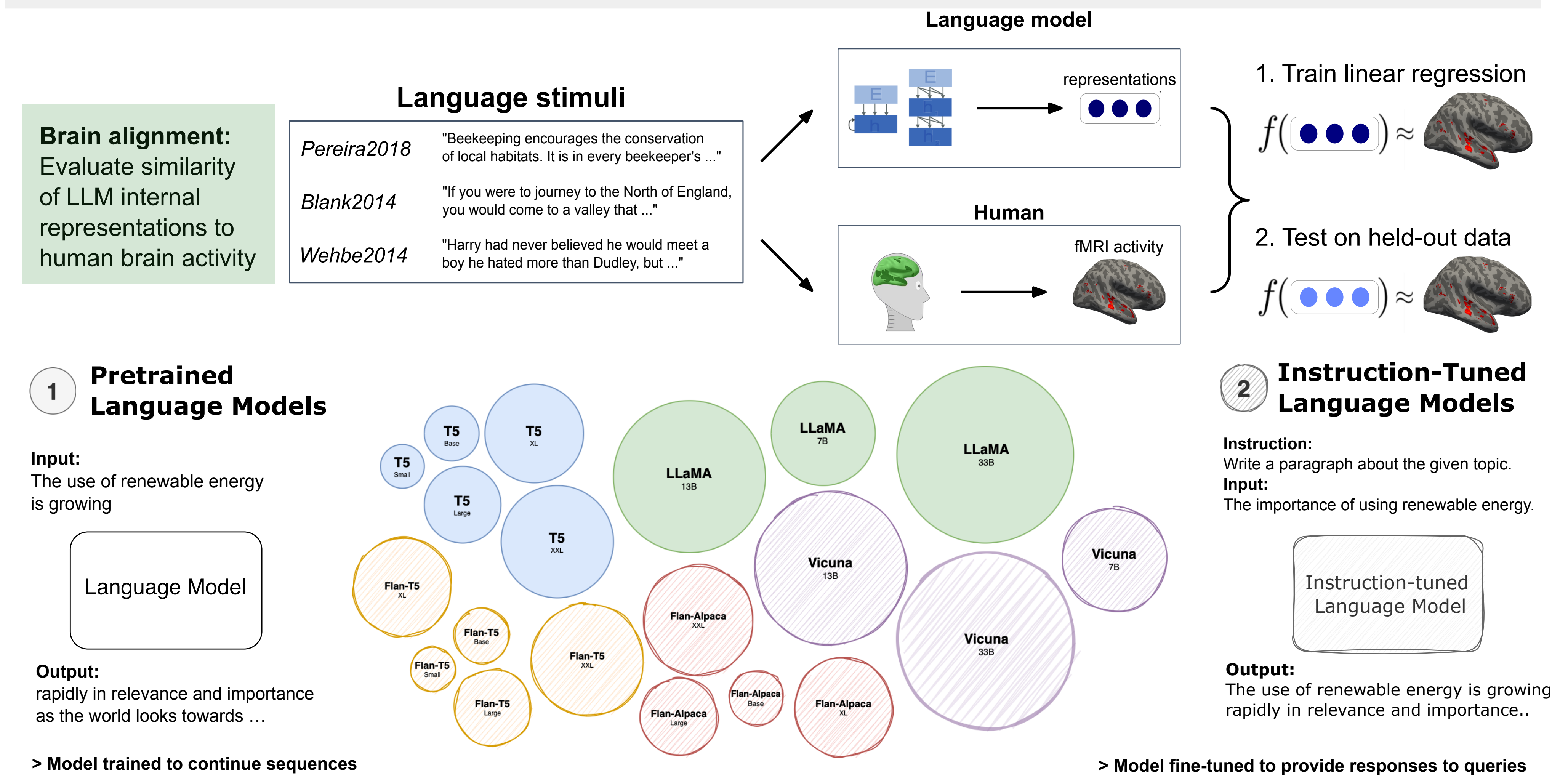
# Instruction-tuning Aligns LLMs to the Human Brain



Khai Loong Aw, Syrielle Montariol\*, Badr AlKhamissi\*, Martin Schrimpf<sup>+</sup>, Antoine Bosselut<sup>+</sup>

\*Equal contribution, <sup>+</sup>Equal supervision / senior authors

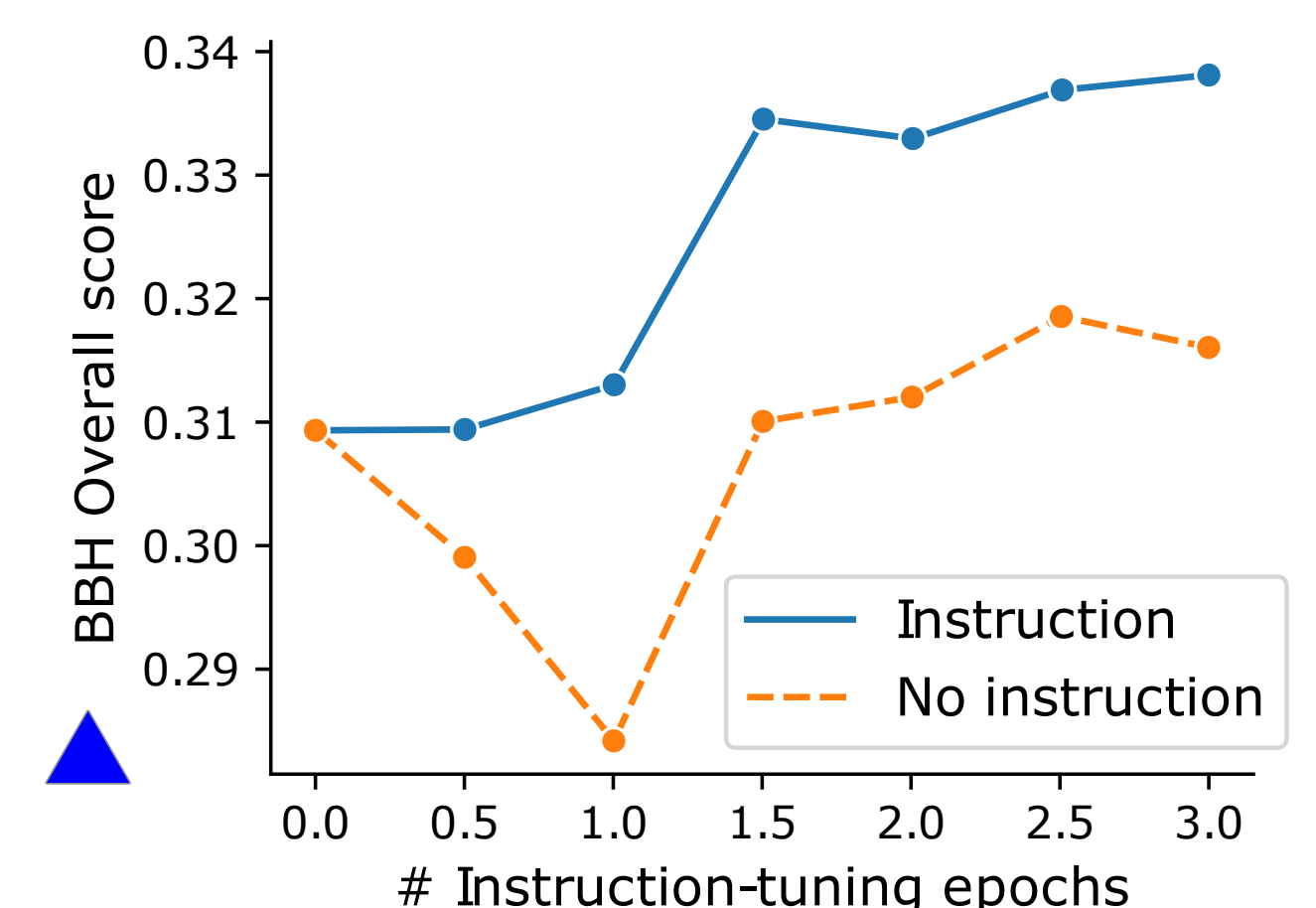
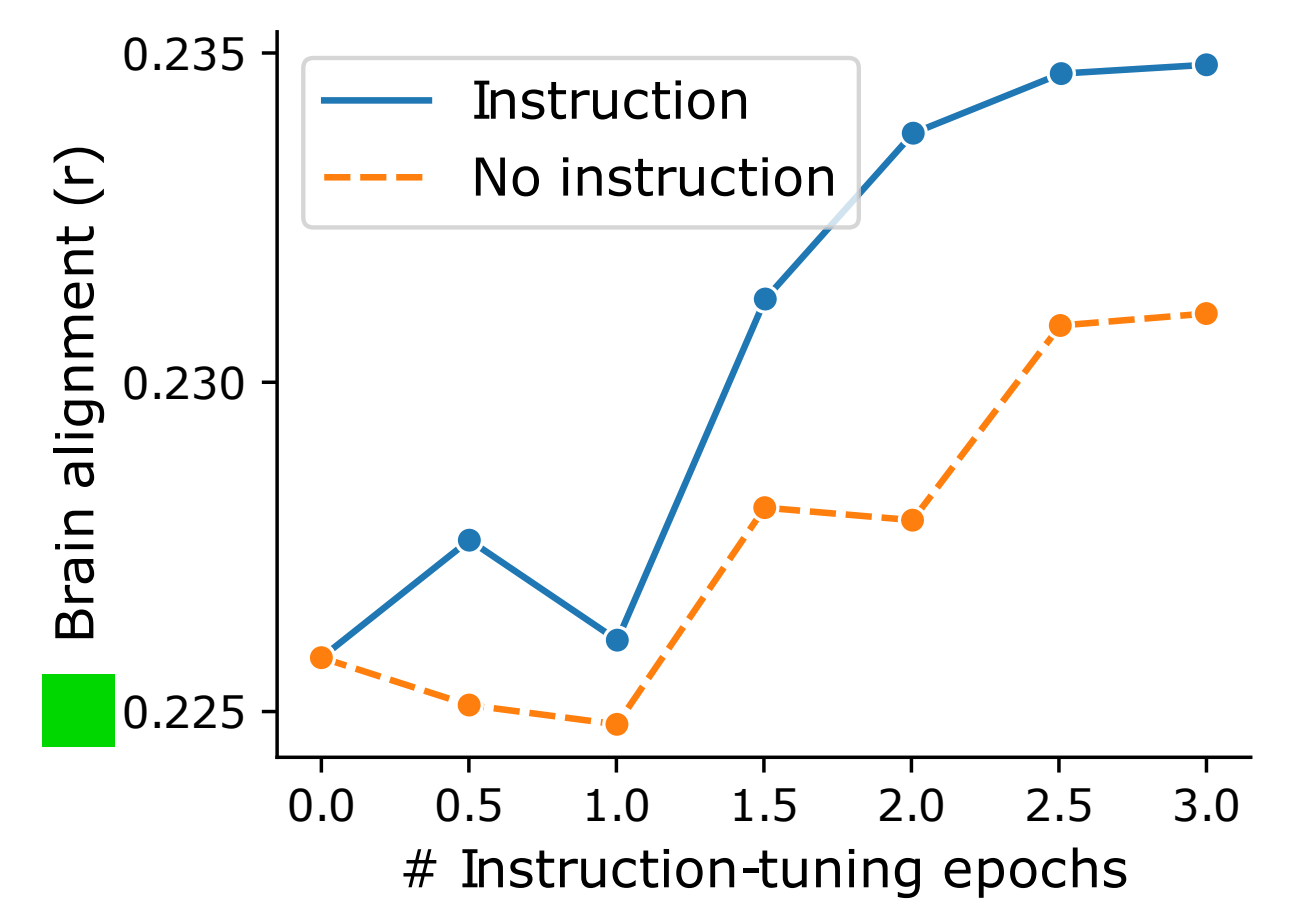
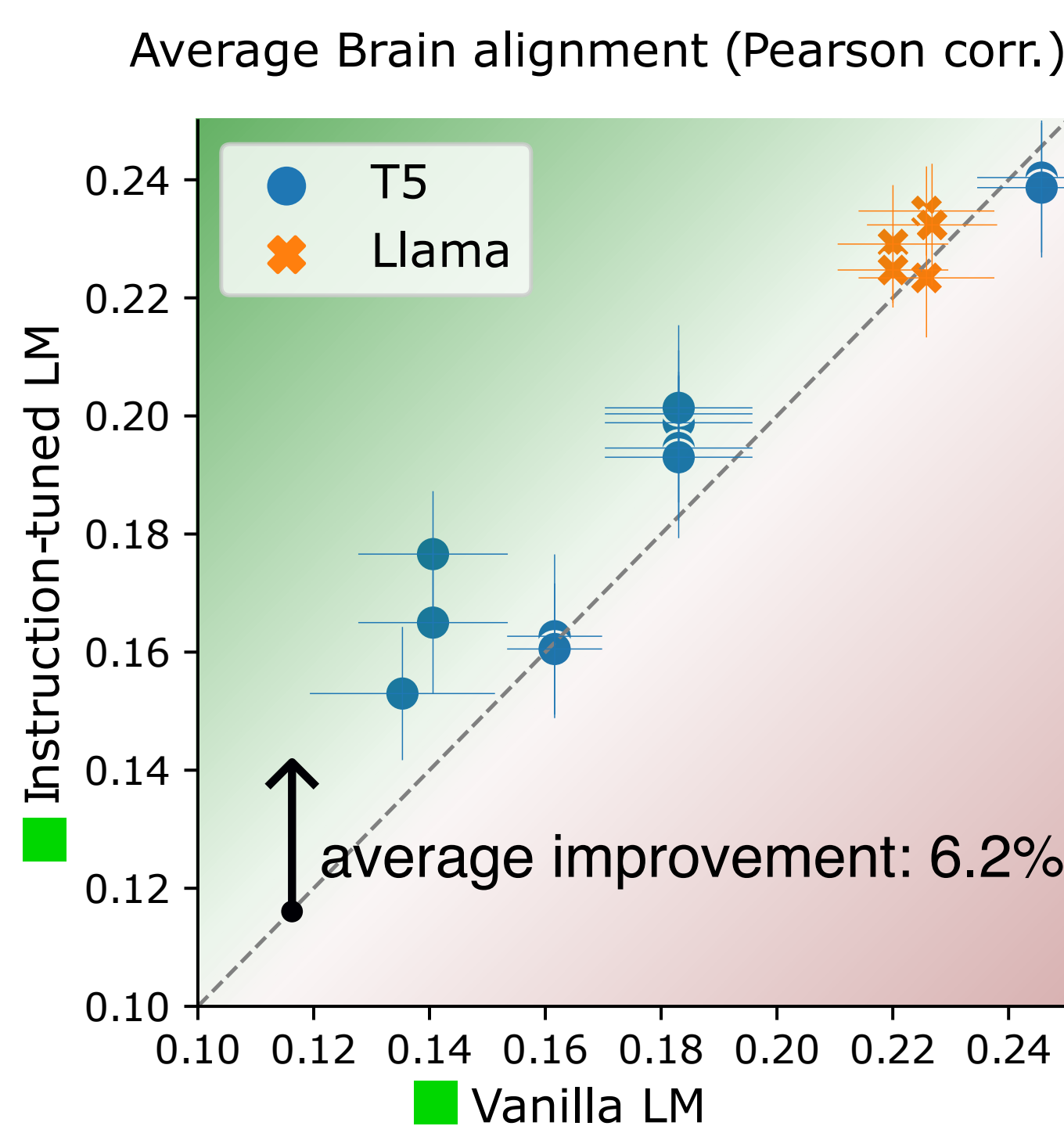
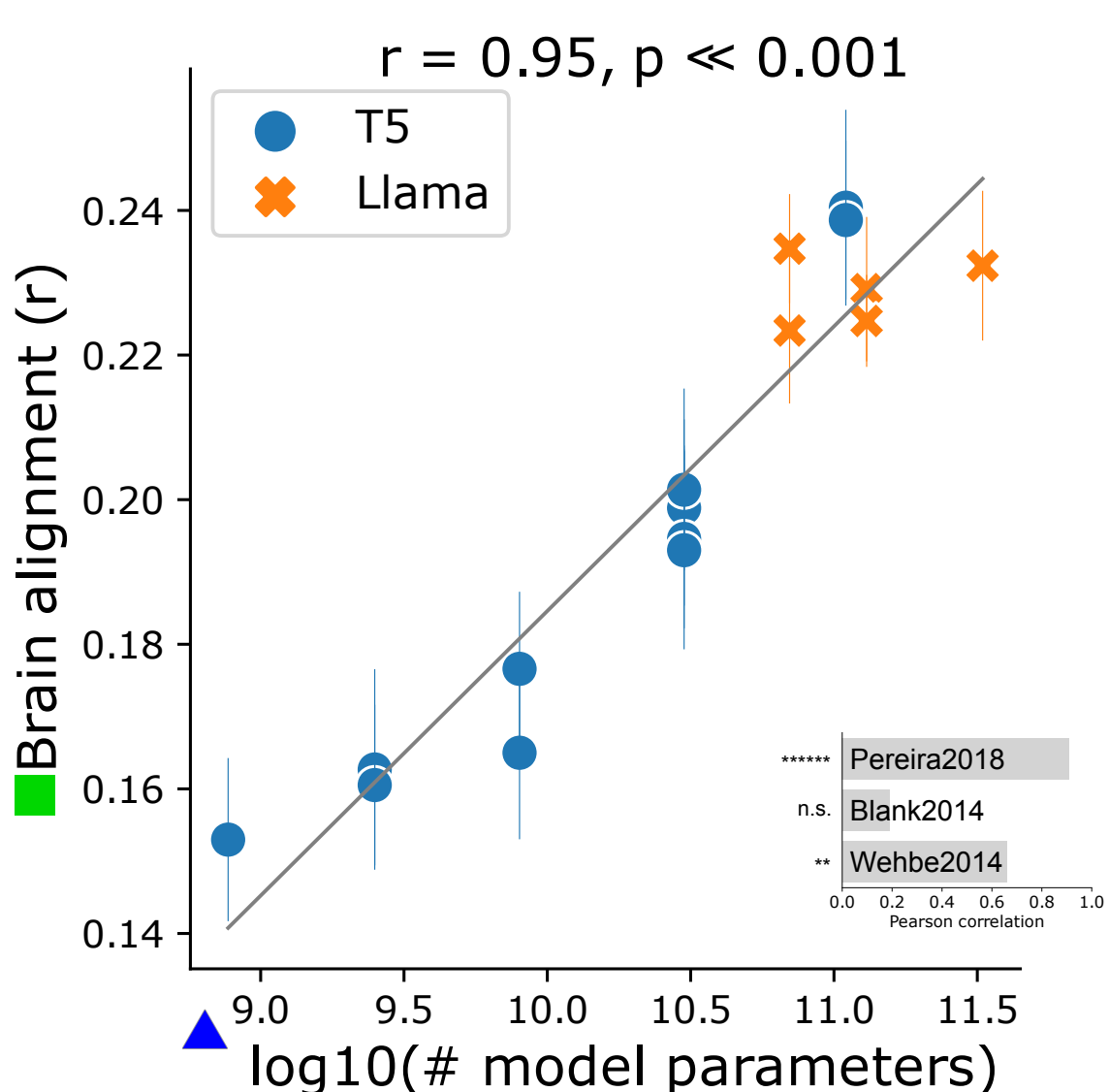
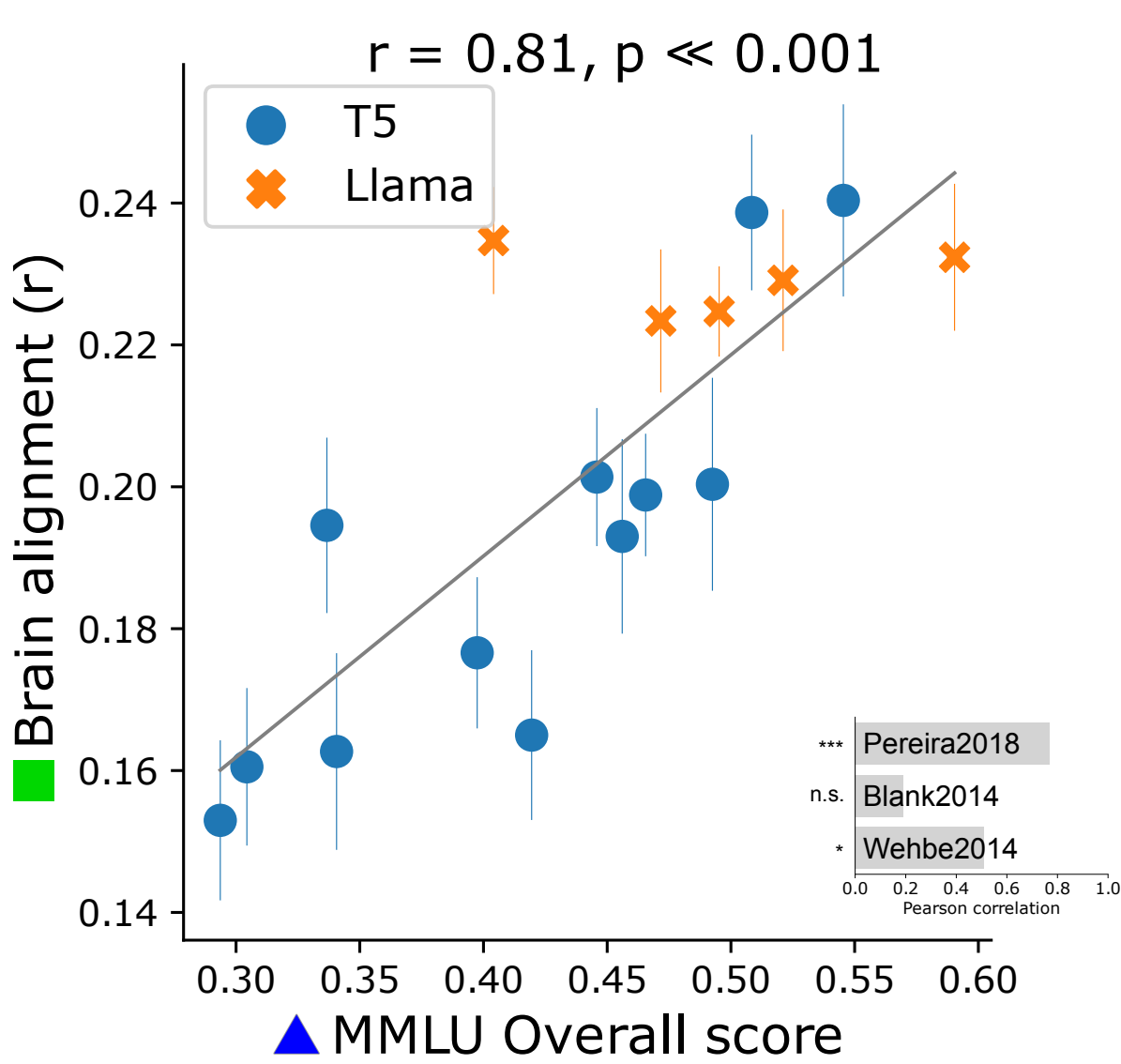
## Preliminaries



## Results

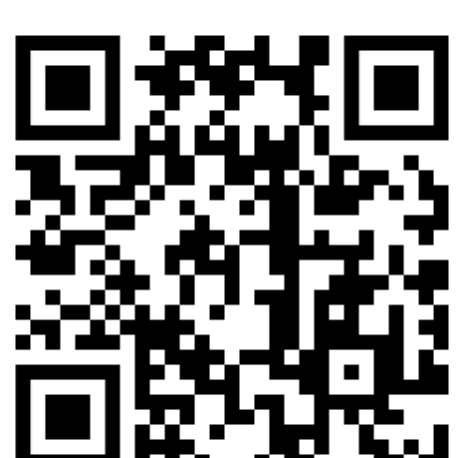
■ Brain alignment ▲ LM properties

**Brain alignment correlates with MMLU & Model Size**      **Instruction tuning improves model brain alignment by 6.2% (avg.)**      **Improvements are due to both training data and process of instruction-tuning**



Correlation between brain alignment and performance **MUCH** higher on **world knowledge** benchmarks

Task category	Brain Alignment Correlation (r)	Corrected p-value	Number of tasks	Average (↑) Performance
MMLU – Overall Score	<b>0.809</b>	<b>0.000329</b>	57	0.36
BBH – Overall score	0.384	0.177	23	0.28
BBH – Algorithmic reasoning	0.194	0.558	8	0.22
BBH – Language understanding	0.163	0.585	3	0.43
BBH – World knowledge	<b>0.679</b>	<b>0.005</b>	5	0.36
BBH – Multilingual reasoning	-0.035	0.895	1	0.19
BBH – Others	0.478	0.083	6	0.27



← Paper Link

@Khai\_Loong\_Aw  
@bkhamssi